

基于多层语义与拓扑融合的异质图方法提升药物-靶标相互作用预测性能

陈紫豪¹, 郭延哺^{1,2}, 宋胜利¹, 郭全明¹, 周冬明^{3,4}

¹郑州轻工业大学软件学院, 河南 郑州 450001; ²东南大学江苏省网络群体智能重点实验室, 江苏 南京 211189; ³湖南信息学院电子科学与工程学院, 湖南 长沙 410151; ⁴云南大学信息学院, 云南 昆明 650500

摘要:目的 为解决药物-靶标相互作用预测中存在的高阶语义依赖建模不足、语义路径融合缺乏自适应性及节点特征过平滑等问题, 提出一种基于多层语义与拓扑融合的异质图预测方法。方法 构建包含药物、蛋白质、副作用、疾病等多类实体的异质图网络, 利用图嵌入技术获取低维特征表示。通过自适应元路径搜索模块, 自动挖掘语义路径组合, 引导高阶语义信息的传播; 构建融合多头注意力的语义聚合机制, 根据上下文信息自动学习各语义路径的重要性, 实现路径间信息的差异化聚合与动态融合; 引入结构感知的门控图卷积模块, 调控特征传播强度, 有效抑制冗余信息, 缓解过平滑问题。最终通过内积操作预测药物与靶标之间的相互作用关系。结果 本文所提方法在公开数据集上, 接收机工作特征曲线下面积(AUC)和精确召回率曲线下面积(AUPRC)分别比现有药物靶标相互作用预测方法的平均性能提高了3.4%和2.4%、3.0%和3.8%。结论 本文设计的药物-靶标相互作用预测方法可有效提取异质生物网络中复杂的高阶语义和拓扑信息, 提升药物-靶标相互作用预测的准确性和稳定性, 可为药物靶标的精准发现和复杂疾病的精准治疗提供技术支撑和理论依据。

关键词: 药物-靶标相互作用; 异质网络; 门控机制; 多头注意力机制; 图卷积网络

A heterogeneous graph method integrating multi-layer semantics and topological information for improving drug-target interaction prediction

CHEN Zihao¹, GUO Yanbu^{1,2}, SONG Shengli¹, GUO Quanming¹, ZHOU Dongming^{3,4}

¹College of Software, Zhengzhou University of Light Industry, Zhengzhou 450001, China; ²Jiangsu Provincial Key Laboratory of Networked Collective Intelligence, Southeast University, Nanjing 211189, China; ³School of Electronic Science and Engineering, Hunan University of Information Technology, Changsha 410151, China; ⁴School of Information Science and Engineering, Yunnan University, Kunming 650500, China

Abstract: Objective To develop a heterogeneous graph prediction method based on the fusion of multi-layer semantics and topological information for addressing the challenges in drug-target interaction prediction, including insufficient modeling of high-order semantic dependencies, lack of adaptive fusion of semantic paths, and over-smoothing of node features. **Methods** A heterogeneous graph network with multiple types of entities such as drugs, proteins, side effects, and diseases was constructed, and graph embedding techniques were used to obtain low-dimensional feature representations. An adaptive metapath search module was introduced to automatically discover semantic path combinations for guiding the propagation of high-order semantic information. A semantic aggregation mechanism integrating multi-head attention was designed to automatically learn the importance of each semantic path based on contextual information and achieve differentiated aggregation and dynamic fusion among paths. A structure-aware gated graph convolutional module was then incorporated to regulate the feature propagation intensity for suppressing redundant information and reducing over-smoothing. Finally, the potential interactions between drugs and targets were predicted through an inner product operation. **Results** Compared with existing drug-target interaction prediction methods, the proposed method achieved an average improvement of 3.4% and 2.4%, 3.0% and 3.8% in terms of the area under the receiver operating characteristic curve (AUC) and the area under the precision-recall curve (AUPRC) on public datasets, respectively. **Conclusion** The drug-target interaction prediction method developed in this study can effectively extract complex high-order semantic and topological information from heterogeneous biological networks, thereby improving the accuracy and stability of drug-target interaction prediction. This method provides technical support and theoretical foundation for precise drug target discovery and targeted treatment of complex diseases.

Keywords: drug-target interaction; heterogeneous networks; gated mechanism; multi-head attention mechanism; graph convolutional networks

收稿日期: 2025-05-21

基金项目: 国家自然科学基金(62403437, 62066047); 河南省科技攻关项目(242102211039); 郑州轻工业大学青年骨干教师培养资助项目(13502010009)

Supported by National Natural Science Foundation of China (62403437, 62066047)

作者简介: 陈紫豪, 在读硕士研究生, E-mail: chenzh000523@163.com

通信作者: 宋胜利, 博士, 教授, 硕士生导师, E-mail: 2003010@zzuli.edu.cn; 周冬明, 博士, 教授, 博士生导师, E-mail: zhoudm@ynu.edu.cn

药物与靶标相互作用(DTI)的预测研究可促进人们对药物耐药性、副作用的理解, 在医学和药物研究领域具有重要意义^[1]。尽管传统的生物化学实验可用于精准预测药物-靶标相互作用, 但十分耗时且昂贵, 其花费取决于实验规模和复杂性^[2]。因此, 设计高效准确的计算方法研究成为当前研究热点。药物-靶标相互作用预测主要包括2个任务: 一是药物-靶标相互作用关系预

测,通常将其视为二元分类问题,用于判断药物和靶标之间是否存在相互作用;二是药物-靶标结合亲和力预测,关注药物与靶标的结合强度^[3]。本文的研究重点主要聚焦于药物-靶标相互作用关系预测任务。

现有的主流药物-靶标相互作用关系预测方法主要是基于机器学习的方法,包括基于特征和基于网络的方法。在基于特征的方法中,利用不同的特征提取策略来提取药物和靶标的生物特征^[4]。例如,Ru等^[5]采用基于距离的Top-n-gram算法结合化合物描述符提取关键特征;Ding等^[6]设计了多视图图正则化传播模型;Song等^[7]融合多维特征并引入交互注意力机制以提升模型性能。然而,这类基于特征的方法对特征定义依赖较强,面对缺失或低质量特征时效果不佳,限制了其实际应用价值。相比之下,基于网络的方法直接建模药物、蛋白质与疾病等多类实体构成的生物网络,通过挖掘节点关系实现全局建模。例如,利用随机游走和扩散分析获取网络特征^[8];融合多源异质信息改进表示学习^[9];构建图生成模型结合因果推断增强预测能力^[10]。这类方法提高了模型对复杂交互的感知能力,但在建模异质网络中的多层语义关系时仍存在局限。

针对这一问题,元路径技术逐渐成为异质网络建模中的关键工具。在生物异质网络中,元路径能够有效捕捉网络的语义信息,全面涵盖代谢过程或生物学原理,为DTI预测的理解和解释提供重要支持^[11]。有研究通过自动提取并加权重要的元路径,增强了模型对药物和靶标关系的建模能力^[12]。有研究提出自适应元图以提升语义聚合效果^[13]。然而,现有方法对于复杂网络的动态变化适应性较弱,且容易受到噪声节点或边的影响。在此背景下,图卷积网络(GCN)^[14]作为一种无需手工特征设计的深度学习方法,在DTI建模中得到广泛应用。Wan等^[15]利用GCN多次传递与聚合图中信息,使节点特征能够体现其拓扑结构关系。有研究通过GCN提取网络的拓扑特征,并结合深度神经网络捕获复杂的非线性关系^[16]。有研究基于异质图表示学习,构建药物-靶标的异质网络,结合不同类型的生物信息,通过端到端的方式优化药物与靶标的表示^[17]。Zhu等^[18]通过构建药物-靶点对网络,利用图结构学习药物和靶标的潜在表示,并结合自监督学习提升药物与靶标相互作用预测的准确性。

然而,现有方法在建模异构网络中的复杂交互关系时,面临以下关键挑战:对高阶语义依赖关系的建模能力有限,多数方法仅关注节点的直接邻居或预设路径,难以有效挖掘药物与靶标之间潜在的深层次语义关联。多语义路径信息融合缺乏自适应性,现有方法大多采用静态或均值策略对语义路径进行处理,未能针对具体语境动态调整各路径的重要性,这种非自适应的融合方式容易掩盖关键语义特征,甚至放大噪声干扰,限制了模

型对复杂网络结构和上下文信息的敏感性与表达能力。节点特征冗余与过平滑问题突出,在多层信息传播过程中,节点特征易被过度平均,噪声逐层积累,导致节点表示趋于同质化,削弱模型的判别性能。

针对现有预测方法的不足,本文提出一种基于多层语义与拓扑融合的异质图预测方法(GMADTI),具体包括3方面创新设计:首先,引入自适应元路径搜索机制,用于自动挖掘关键语义路径组合,引导高阶语义信息的传播与聚合,从而提升模型对深层语义依赖的建模能力。其次,设计融合多头注意力的异质语义聚合模块,通过多个注意力子空间并行建模语义路径间的关联性,自适应地分配不同语义路径的权重。该方法有效提升了语义信息融合的自适应性,使模型能够根据上下文结构自动调整不同路径的贡献程度,从而增强对复杂交互关系的表达能力。第三,设计基于门控机制的结构感知图卷积模块,动态调控信息传播强度,有效抑制冗余与噪声,缓解过平滑现象。最终,基于药物与靶标的嵌入向量,通过内积操作预测其潜在相互作用关系。相较于现有DTI预测方法,GMADTI在高阶语义挖掘、自适应语义融合与特征传播控制3方面均实现了结构级优化,切实提升了模型对复杂异构交互关系的建模能力。

1 基于多层语义与拓扑融合的异质图药物-靶标相互作用预测算法设计

本文提出的GMADTI模型的整体框架如图1。该预测算法主要包含4个部分:(1)嵌入层模块:构建了一个包含药物、蛋白质、副作用、疾病的异质网络,并采用Node2Vec^[19]算法对药物节点和蛋白质节点进行低维嵌入表示,提取全局拓扑信息。(2)自适应搜索元路径模块:在异质网络中为药物和蛋白质节点自适应搜索元路径,挖掘潜在的语义关联。(3)门控多头注意力模块:引入了门控机制,能够动态地学习异质网络中的特征,并且利用多头注意力机制深入挖掘药物和蛋白质之间复杂的交互模式,从而生成药物和靶标的拓扑信息表示。(4)预测器模块:将药物和蛋白质的嵌入向量进行内积运算,以预测潜在的药物-靶标相互作用。

1.1 节点编码模块

Node2Vec算法是一种常见的图数据预处理方法,本文采用Node2Vec算法来对异质网络中的节点进行编码。首先节点编码模块通过执行多次随机行走,生成一系列节点序列,反映节点在图中的连接模式及其邻域关系。然后,利用Word2Vec^[20]中的Skip-gram模型对这些序列进行训练,捕捉局部邻域关系和全局结构信息,获得节点的低维嵌入向量。

1.2 自适应元路径模块

1.2.1 自适应元路径的定义 为了在药物靶标相互作用

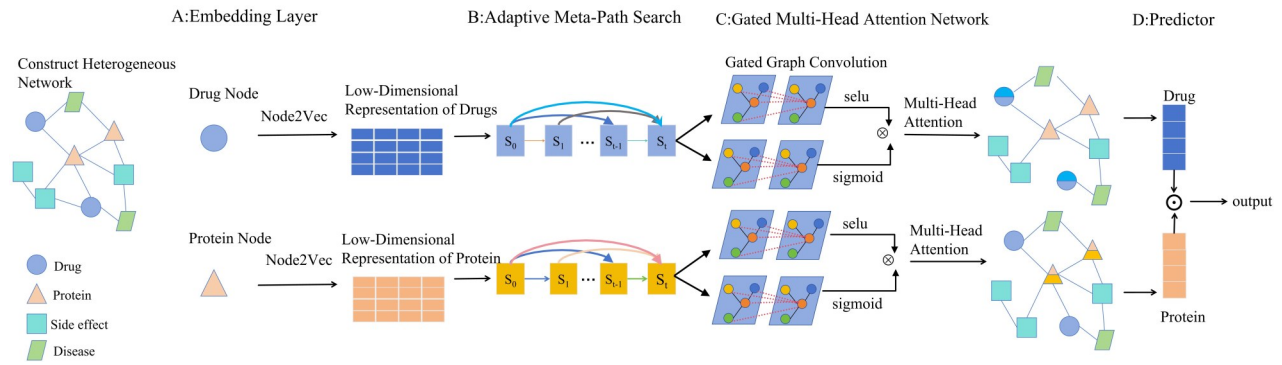


图1 GMADTI总体框架示意图

Fig.1 Overall framework of GMADTI.

预测中有效识别元路径^[13],本文设计了一个自适应元路径模块。该模块的核心在于构建自适应元图,这是一个有向无环图,记作 $M=(V_m, E_m)$ 。其中,节点集 V_m 表示异质网络中每次信息传播后的节点特征集合 $\{S_0, S_1, \dots, S_T\}$, S_i 是第*i*次传播后的特征状态,节点的数量由异质网络的信息传播次数决定。边类型集合 E_m 表示所有可能的信息传播方式。例如,若从 S_0 到 S_1 的定向链接标记为“蛋白质→疾病”,则这意味着 S_1 是通过将 S_0 中的“蛋白质”节点的特征聚合到“疾病”节点的特征来获得的。在所提出的自适应元图结构中,任意前一状态 $S_i \in \{S_0, S_1, \dots, S_{T-1}\}$ 均可能通过特定传播方式影响当前状态 S_T ,从而在异质网络的不同状态之间产生跳跃结构,从而在不同传播状态之间构建跳跃连接。这使得模型能够更充分地捕捉异质网络中的复杂语义信息,这是自适应元图的第1个特点。

自适应元图的另1个特点是元图中的每条链路都是自适应确定的,即异质网络的前一状态是否影响当前状态,自适应地确定影响当前状态的方式,从而将异质网络中的所有边类型作为可能的信息传播模式。此外,引入了两个辅助传播方式: L_1 表示当前状态等于前一状态, L_2 表示前一状态无法影响当前状态。最终,在自适应元图中,节点之间有12种可能的连接: $\{L_{DP}, L_{PD}, L_{DS}, L_{SD}, L_{DE}, L_{ED}, L_{PE}, L_{EP}, L_{DD}, L_{PP}, L_1, L_2\}$,其中前10种连接类型对应异质网络中的实际边类型(如 L_{DP} :药物→蛋白质、 L_{ED} :疾病→药物等),最后两类为新增设计。

1.2.2 自适应元路径的构造 对于DTI预测,应用自适应元图来引导异质网络中的信息聚合,以获得药物和蛋白质的特征。对于构建自适应元图的方法,首先,自适应元图中节点的数目取决于信息在异质网络中传播的次数,假设节点的特征在异质网络中传播*T*次,则自适应元图中的节点为 $\{S_0, S_1, \dots, S_T\}$ 。在自适应元图中,节点对之间的可能连接方式依据其状态在信息传播过程中的相对位置进行设定。对于给定的2个节点 S_i 和 S_t ($0 \leq i, t \leq T, i \in \mathbb{N}$ 且 $t \in \mathbb{N}$),需判断 S_i 是否为 S_t 的前一状态,

以及 S_t 是否为最终状态。当 $i=t-1$ 且 $t < T$ 时,由于第*i*次信息传播中的节点特征会以某种方式影响第*t*次传播,因此从 S_i 到 S_t 的可能连接是集合 E_m 中除 L_2 外的所有类型。当 $i < t-1$ 且 $t < T$ 时,此时状态 S_i 可能不会影响状态 S_t ,因此可能得连接包括 L_2 ,当 $t=T$ 时,即最后一层状态,那么从 S_{t-1} 到 S_t 的可能连接进一步限制为与药物或靶点相关的连接。

一方面,本文选择与约束条件(R)相关的可能连接作为“*→药物”形式的类型。另一方面,选择与蛋白质相关的“*→蛋白质”的类型,以更新蛋白质的节点特征。例如,在异质网络中,选择4种模式来更新药物的节点特征,即: $R = \{L_{PD}, L_{SD}, L_{ED}, L_{DD}\}$,分别对应蛋白质→药物,副作用→药物,疾病→药物,药物→药物。若 $i=t-1$ 且 $t=T$,则仅保留上述目标相关边;若 $i < t-1$ 且 $t=T$,则在此基础上引入 L_1 和 L_2 以增强灵活性。综上节点 S_i 到 S_t 的可能连接方式可归纳为如下公式:

$$C_{t,i} = \begin{cases} E_m - \{L_2\}, & i = t - 1, t < T \\ E_m, & i < t - 1, t < T \\ R, & i = t - 1, t = T \\ R \cup L_1 \cup L_2, & i < t - 1, t = T \end{cases} \quad (1)$$

然后,从 S_i 到 S_t 的连接类型将从所有可能连接中自适应地选择。在所提出的GMADTI模型中,给定从 S_i 到 S_t 的每一种可能连接,分配一个对应的可学习参数 $\theta_{t,i}^n$,用于表示该连接被选中的可能性。例如,对于从 S^0 到 S^2 的连接“蛋白质→药物”(记作 L_{PD})的可能性被分配一个参数 $\theta_{0,2}^1$ 。在所有连接中,具有最大参数值的连接(即 $\theta_{t,i}^n = \max(\theta_{t,i}^0, \dots, \theta_{t,i}^{11})$)将更可能被选为最终路径。此外,引入了随机采样策略,提高元图结构的多样性。从候选连接集合 $C_{t,i}$ 中以概率 p_i 随机选择一条边;以概率 $1-p_i$ 选择参数值最大的连接。该参数 $p_i \in (0, 1)$ 被设为一个较小的值,用于在训练初期促进不同信息传播路径的探索,并在传播轮数增加时逐渐减小,最终趋近于0,以增强结构的确定性。最终,从 S_i 到 S_t 的连接类型确定方式如下所示:

$$C_{t,i}^m = \begin{cases} \theta_{t,i}^m, & 1 - p_i \\ \text{rand}(C_{t,i}), & p_i \end{cases} \quad (2)$$

其中 $\text{rand}(C_{t,i})$ 表示从连接集合 $C_{t,i}$ 中进行均匀随机采样, m 表示当前具有最大连接概率的边类型索引。 $\theta_{t,i}^m$ 的计算是自适应选择链路类型的关键, 这里采用基于网络结构搜索的方法(DiffMG^[21])用于测量的值。

1.3 门控图卷积特征提取模块

为了增强图神经网络在复杂结构数据上的特征提取能力, 本文设计了一个门控图卷积特征提取模块。门控机制的引入有效提升了模型对图中不同节点特征的选择能力, 抑制了无关信息的干扰, 从而提高了模型的表达能力和鲁棒性。该特征提取模块在传统图卷积操作中引入门控机制^[22], 使信息流的传递不仅依赖于图的拓扑结构, 还能通过学习的门控权重自适应地调整特征信息的传播强度。首先, 该特征提取采用标准图卷积层对输入特征 H^k 进行处理, 通过邻域聚合操作更新每个节点的特征表示。如公式(3)所示:

$$H_1^{k+1} = \tau(D^{\frac{1}{2}} L_k D^{\frac{1}{2}} H^k W^k) \quad (3)$$

其中 L^k 是通过自适应元路径学习得到的邻接矩阵; $D^{\frac{1}{2}}$ 是 L^k 的归一化的度矩阵; H^k 是第 k 层的节点表示; W^k 是可学习的权重矩阵; τ 是 SELU^[23] 激活函数。

接着, 特征提取引入一个独立的图卷积层作为门控结构。该层同样对输入特征 H^k 进行卷积操作, 但其输出经过 Sigmoid 激活函数^[24] 处理, 将值限制在 0~1。门控层生成的输出可视为动态的“权重”或“开关”, 用于调控每个节点的特征流动。通过 Sigmoid 激活函数, 门控机制能够灵活地调整特征信息的通过比例。在这一过程中, 关键信息得以保留, 而不相关或噪声信息则被有效抑制。如公式(4)所示:

$$g_k = \sigma(D^{\frac{1}{2}} L_k D^{\frac{1}{2}} H^k W^t) \quad (4)$$

其中 W^t 是另一个可学习的权重矩阵, σ 是 Sigmoid 激活函数。最后, 门控层的输出与图卷积层的输出逐元素相乘, 实现了对特征流的控制。如公式(5)所示:

$$H_2^{k+1} = g_k \otimes H_1^{k+1} \quad (5)$$

其中 \otimes 表示逐元素相乘操作, g_k 是门控层输出的门控值, 动态调整的“权重”或“开关”, 控制每个节点的特征流动。经过门控处理的输出 H_2^{k+1} 通过批标准化进一步调整, 以增强模型训练的稳定性。批标准化不仅能够加速收敛, 还能有效缓解梯度消失或爆炸等问题, 从而降低模型过拟合的风险。最终的输出是经过门控筛选的节点特征。这些特征不仅包含了通过图结构聚合的关键信息, 还通过动态过滤有效地剔除了无关或噪声信息, 确保了特征表达的质量和模型的鲁棒性。如公式(6)所示:

$$H^{k+1} = \frac{H_2^{k+1} - \mu}{\sqrt{\sigma^2 + \varepsilon}} \times \gamma + \beta \quad (6)$$

其中 μ 和 σ^2 是当前批次的均值和方差, γ 和 β 是可学习的缩放和偏移参数, ε 是为了数值稳定性添加的一个很小的常数, 防止除以零。

1.4 多头注意力模块

为了有效捕捉节点间复杂的交互模式, 本文设计了一个多头注意力^[25] 模块。该模块通过多个注意力头并行计算每个节点的重要性权重, 从而在不同的表示空间内学习节点间多样化的关系特征。首先, 通过 L_1 层, 每个注意力头将输入节点的特征映射到隐藏层维度, 以便更好地捕捉节点间的复杂关系。接着, L_2 层将变换的特征进一步映射为一个标量值, 用于表示每个节点在该头下的注意力得分。这些得分反映了节点在信息聚合过程中的重要性, 进而影响最终的特征聚合。通过多个独立的注意力头, 模型能够在不同的子空间中并行学习不同的关系模式, 增强了模型在复杂异质网络中的表达能力和捕捉潜在关系的能力。多头注意力模块特征计算方法如公式(7):

$$\text{attn}_i^k = L_2(\tanh(L_1(h^k))) \quad (7)$$

其中, h^k 表示输入的特征表示, L_1 和 L_2 分别是两层线性变换函数, \tanh 是激活函数。然后将各头的输出拼接在一起, 拼接结果通过 softmax^[26] 进行归一化操作, 生成最终的节点特征嵌入。如公式(8)所示:

$$\text{attn}^k = \text{softmax}(f_{\text{concat}}(\text{attn}_1^k, \text{attn}_2^k, \text{attn}_3^k, \text{attn}_4^k)) \quad (8)$$

其中 $\text{attn}_1^k, \text{attn}_2^k, \text{attn}_3^k, \text{attn}_4^k$ 分别表示第 k 层中每个注意力头的输出, f_{concat} 表示将它们沿最后一个维度进行拼接, softmax 是归一化函数。最终, 模型利用这些权重对隐藏状态进行加权求和, 从而生成包含多维度信息的输出。

1.5 损失函数模块

为优化模型分类性能与嵌入表示, 本文提出了一种结合二元交叉熵损失函数(BCE Loss)^[27] 和负对数似然损失函数(NLL Loss)^[28] 的“双重损失函数”策略。具体而言, BCE 损失通过 Sigmoid 函数实现概率映射, 提供稳定的优化过程, 其样本级梯度计算有效缓解了 DTI 任务中常见的正负样本不平衡问题, 并通过全局监督信号确保模型的整体优化方向。而 NLL 损失基于对比学习机制, 通过最大化正样本对似然和最小化负样本对似然, 增强特征表示的判别性: 一方面显式调整嵌入空间结构, 使相互作用药物-靶标对聚集而非相互作用对分离; 另一方面通过 Log-Sigmoid^[29] 运算强化对边界样本的学习, 提升模型鲁棒性。

二者的优势互补形成了多粒度监督机制: BCE 损失在样本层面提供分类指导, NLL 损失在特征空间建立关联约束, 这种局部与全局相结合的学习策略使模型能

够更全面地捕捉药物-靶标相互作用的潜在规律,从而获得更准确可靠的预测结果。本文定义的双重损失函

数具体形式如公式(9)所示:

$$\eta = \eta_{BCE} + \eta_{NLL}$$

$$\eta_{BCE} = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(\sigma(s_i)) + (1 - y_i) \cdot \log(1 - \sigma(s_i))]$$

$$\eta_{NLL} = -\frac{1}{N} \sum_{i=1}^N \left[\log(\sigma(f(out_s^p \cdot out_t^p))) + \log(1 - \sigma(f(out_s^n \cdot out_t^n))) \right]$$
(9)

其中 η_{BCE} 为二元交叉熵损失函数, η_{NLL} 为负对数似然损失函数, y_i 为样本标签, s_i 为模型输出的原始值, σ 表示 Sigmoid 函数, out_s^p 和 out_t^p 分别为药物输出的正样本和负样本, out_s^n 和 out_t^n 分别为靶标输出的正样本和负样本, $f(out_s^p \cdot out_t^p)$ 是正样本对的相似度, $f(out_s^n \cdot out_t^n)$ 是负样本对的相似度。

2 结果

2.1 数据集介绍

为了测试 GMADTI 的性能表现, 本文选用 Luo^[8] 和 Zheng^[30] 提供的 2 组公开数据集进行实验验证。Luo 数据集涵盖了 4 大主要实体类别: 药物、蛋白质、疾病以及副作用。这些实体之间的相互作用关系包括药物-靶标、疾病-蛋白质以及药物-副作用等多种类型, 数据的多样性为研究复杂的生物网络提供了基础, 数据集中节点和边的信息如表 1。相比之下, Zheng 数据集不仅包含药物和靶标的基本信息, 还引入了丰富的药物和蛋白质属性特征, 例如药物的取代基信息、化学结构描述、蛋白质的基因本体注释以及蛋白质的氨基酸序列特征等, 数据集中节点和边的信息如表 2。这些附加的高维特征有助于模型更全面地捕捉药物和靶标之间的潜在相互作用模式, 为 DTI 预测提供更细粒度的信息支持。

表 1 Luo 数据集中节点和边的信息

Tab.1 Information of the nodes and edges in the Luo dataset

Node type	Num	Edge type	Num
Drug	708	Drug-drug (interaction)	10 036
Protein	1512	Drug-drug (interaction)	501 264
Disease	5603	Drug-protein	1923
Side effect	4192	Drug-disease	199 214
		Drug-side effect	80 164
		Protein-disease	1 596 745
		Protein-protein (interaction)	7363
		Protein-protein (similarity)	2 286 144

2.2 实验设置和评估指标

本文基于 PyTorch 框架实现了所提出的 GMADTI 模型, 并采用 NAdam^[31] 优化器进行训练。优化器参数

表 2 Zheng 数据集中节点和边的信息

Tab.2 Information of the nodes and edges in the Zheng dataset

Node type	Num	Edge type	Num
Drug	1094	Drug-drug	1 196 836
Protein	1556	Drug-drug	11 819
Chemical structure	881	Drug-chemical substructure	133 880
Side effect	4063	Drug-side effect	122 792
Substituent	738	Drug-side effect	20 798
GO term	4098	Protein-GO term	35 980
		Protein-protein	2 421 136

设置如下: 学习率设定为 $2e-3$, 权重衰减率设为 0, 注意力头数为 4, 训练轮数为 100。Node2Vec 的超参数设置为: 游走长度为 100, 行走次数为 10, 前向和后向跳跃因子均设为 1。在 Luo 和 Zheng 数据集上进行实验时, 唯一不同的参数是隐藏层的维度: Luo 数据集采用 64 维, Zheng 数据集则为 256 维, 其他参数设置均保持一致。各基线方法的参数均严格按照原始文献中的推荐设置进行配置, 确保比较的公平性与有效性。

为全面评估 GMADTI 模型的性能, 本文在 2 个基准数据集上采用五折交叉验证进行实验。考虑到未知药物-靶标相互作用的数量远大于已知相互作用, 本文对未知样本进行了欠采样, 使正负样本数量保持一致, 构建平衡的数据集。每轮交叉验证中, 随机选取 60% 的正负样本作为训练集, 20% 作为验证集用于参数调整, 其余 20% 用于测试模型性能。训练集用于模型学习, 验证集用于优化模型参数, 测试集用于性能评估。本文采用了接收机工作特征曲线下面积 (AUC) 和精确召回率曲线下面积 (AUPRC) 作为主要评价标准。实验在 2 个数据集上均进行 5 次五折交叉验证, 以反映其预测能力与稳定性。

2.3 与其它方法比较

为了全面评估提出的 GMADTI 模型的性能, 本文选取一些代表性的药物-靶标相互作用预测方法作为基线模型, 涵盖不同建模策略。其中, 包括基于异质网络信息建模的方法, 例如 DTINet^[8]、NeoDTI^[15]、EEG-DTI^[17]、CE-DTI^[10] 及 GSRF-DTI^[18] 等; 基于图神经网络

结构建模的方法,例如 GCN-DTI^[16]、IMCHGAN^[12]、HampDTI^[32]、SGCL-DTI^[33]和SHGCL-DTI^[34]等;以及关注元路径语义挖掘与结构自适应建模的方法,例如 MIDTI^[7]和AMGDTI^[13]。这些模型均在 Luo 与 Zheng 两个标准数据集上进行对比实验,评估所提出方法在不同建模机制下的整体表现。

汇总 GMADTI 与多种基线方法在 Luo 和 Zheng 数据集上的性能结果(表3、4),相比现有 12 种机器学习方法,GMADTI 的 AUC 和 AUPRC 平均提高了 3.2% 和 3.1%。其中,在 Luo 数据集上,GMADTI 的 AUC 和 AUPRC 均达到 0.987±0.002,比当前表现最优的 AMGDTI 均提升了 1%。在 Zheng 数据集上,GMADTI 的 AUC 和 AUPRC 分别达到 0.987±0.002 和 0.983±0.001,相较于 AMGDTI 分别提升了 1.4% 和 1.2%。

表3 在 Luo 数据集上与基线方法的对比结果

Tab.3 Comparison results with baseline methods on the Luo dataset

Methods	Luo dataset	
	AUC	AUPRC
DTINet	0.879±0.004	0.906±0.003
NeoDTI	0.955±0.003	0.889±0.004
GCN-DTI	0.918±0.005	0.897±0.005
IMCHGAN	0.956±0.004	0.959±0.003
HampDTI	0.928±0.003	0.927±0.005
SGCL-DTI	0.977±0.002	0.976±0.002
GSRF-DTI	0.977±0.002	0.980±0.003
EEG-DTI	0.954±0.003	0.964±0.004
MIDTI	0.978±0.003	0.970±0.002
CE-DTI	0.976±0.002	0.976±0.002
SHGCL-DTI	0.957±0.004	0.958±0.003
AMGDTI	0.977±0.002	0.977±0.002
GMADTI	0.987±0.002	0.987±0.002

2.4 消融实验

本文设计了 5 个消融实验研究门控图卷积、多头注意力机制、双重损失函数以及元路径搜索机制对模型性能的影响。根据各组件的配置差异,构建了 5 种模型变体,分别命名为 GMADTI-GCN、GMADTI-NoGateWeight、GMADTI-SingleHead、GMADTI-SingleLoss 和 GMADTI-FixedPath。GMADTI-GCN 是将门控图卷积替换为标准的图卷积网络,GMADTI-NoGateWeight 在门控图卷积模块中将门控权重恒设为 1,即去除节点特征的动态调节能力,仅保留图卷积主通路;GMADTI-SingleHead 将多头注意力机制替换为单头注意力机制,GMADTI-SingleLoss 是去除双重损失

表4 在 Zheng 数据集上与基线方法的对比结果

Tab.4 Comparison results with baseline methods on the Zheng dataset

Methods	Zheng dataset	
	AUC	AUPRC
DTINet	0.889±0.004	0.900±0.004
NeoDTI	0.946±0.003	0.846±0.005
GCN-DTI	0.922±0.004	0.914±0.004
IMCHGAN	0.946±0.002	0.929±0.003
EEG-DTI	0.942±0.003	0.941±0.003
SGCL-DTI	0.968±0.002	0.968±0.002
SHGCL-DTI	0.957±0.003	0.961±0.003
CE-DTI	0.972±0.003	0.972±0.002
MIDTI	0.954±0.002	0.949±0.004
AMGDTI	0.973±0.004	0.971±0.002
GMADTI	0.987±0.002	0.983±0.001

函数,仅使用二元交叉熵损失函数,GMADTI-FixedPath;移除自动元路径搜索模块,固定使用 Drug→Disease→Protein→Protein 元路径作为结构输入。实验结果显示(图2、3),GMADTI 在 2 个数据集上均优于所有消融变体。

在 Luo 数据集上,去除门控机制(GMADTI-GCN)后,模型的 AUC 下降了 2.1%,AUPRC 下降了 2.4%;去除门控机制中的权重模块(GMADTI-NoGateWeight)后,AUC 和 AUPRC 分别下降了 0.6% 和 0.5%;去除多头注意力机制(GMADTI-SingleHead)后,AUC 和 AUPRC 均下降约 0.3%;移除双重损失函数(GMADTI-SingleLoss)后,AUC 和 AUPRC 分别下降约 0.3% 和 0.4%;替换为固定元路径(GMADTI-FixedPath)后,性能均下降约 0.4%。在 Zheng 数据集上,去除门控机制后,AUC 和 AUPRC 分别下降了 0.9% 和 1.0%;去除门控权重模块和多头注意力机制后,性能均下降约 0.3%;去除双重损失函数后,AUC 和 AUPRC 分别下降 0.2% 和 0.4%;固定元路径代替自适应搜索后,性能下降约 0.2%。

2.5 参数敏感性分析

为评估不同超参数对模型性能的影响,本文在 Luo 数据集上进行了参数敏感性分析,重点考察多头注意力机制中的头数、学习率和嵌入维度等关键参数。通过对这些超参数进行系统调节,以获得最佳的超参数组合,从而优化预测模型的精度和泛化能力。

多头注意力机制中的头数直接决定了模型在多个子空间并行学习和聚合节点特征的能力。在本实验中,模型对注意力头数在 {2,4,6,8} 范围内进行探索。当注意力头数设置为 4 时,模型在 AUC 和 AUPRC 指标上均表现最佳(图4)。学习率是影响模型训练收敛速度和稳定性的关键超参数,为研究学习率对模型性能的影响,

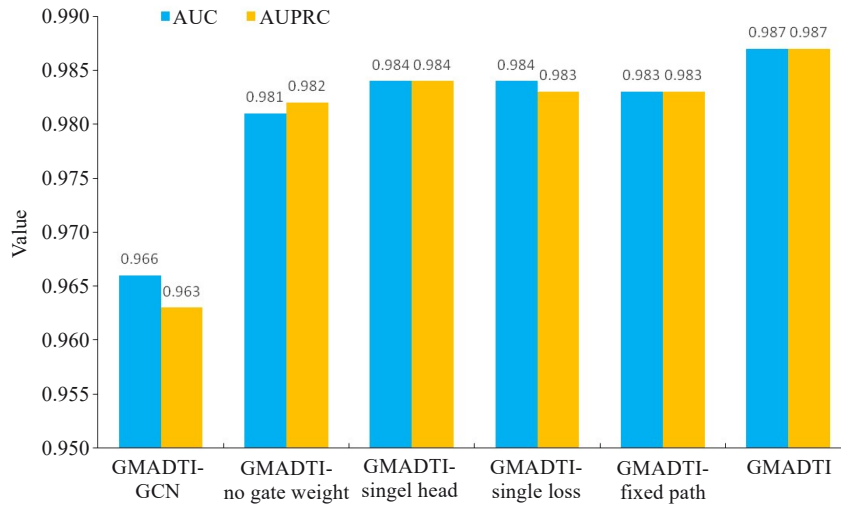


图2 在Luo数据集上GMADTI模型及其变体模型比较

Fig.2 Comparison of GMADTI model and its variants on the Luo dataset.

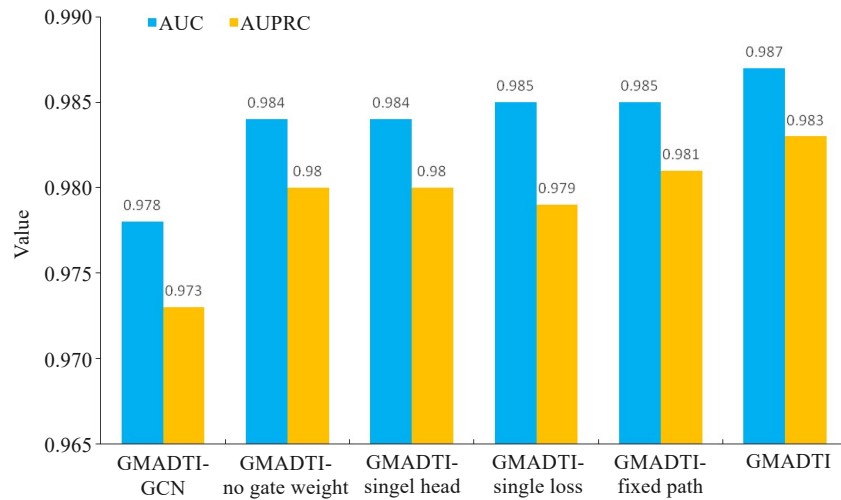


图3 在Zheng数据集上GMADTI模型及其变体模型比较

Fig.3 Comparison of GMADTI model and its variants on the Zheng dataset.

本实验将学习率设置为 $\{5e-4, 1e-3, 2e-3, 5e-3\}$,并观察模型在这些学习率下的训练过程以及最终表现。当学习率设置为 $2e-3$ 时,模型的性能最佳,既能保持较快的收敛速度,又能达到较高的预测精度(图5)。嵌入维度决定了特征表示的细节程度和信息的多样性。本实验选择了4种嵌入维度(32、64、128和256)进行测试,实验结果显示,当嵌入维度设置为64时,模型在Luo数据集上表现最佳,能够在特征表达能力和计算效率之间达到较好的平衡(图6)。

2.6 模型可视化结果

2.6.1 t-SNE 嵌入可视化结果 为验证模型在特征表达方面的有效性,本文以 Luo 数据集为例,基于训练后的药物与靶标嵌入向量,采用 t-SNE 方法对高维特征进行降维,并进行二维可视化分析(图7),图中蓝色与红色点分别对应标签为0和1的节点,代表药物与靶标。药物

与靶标节点在嵌入空间中形成2个主要聚簇,类别边界较为清晰,簇内呈现一定的聚集趋势。部分节点分布于簇边缘,表现为离群点。

2.6.2 注意力权重热图结果 进一步分析模型在信息聚合过程中的关注机制,本文基于 Luo 数据集,引入注意力权重的可视化分析。通过提取训练后模型的注意力分布矩阵,并绘制热力图,展示模型在不同注意力头与传播步骤组合下对样本的关注程度,以分析其语义偏好。药物侧与靶标侧的注意力热力图(图8、9),纵轴表示样本索引(0~9),对应不同的药物-靶标对;横轴表示注意力头与传播步骤的组合(编号0~3)。颜色深浅表示注意力权重大小,深蓝色代表较高权重,浅黄色表示较低权重。

3 讨论

为全面验证 GMADTI 的有效性 with 稳定性,本文从

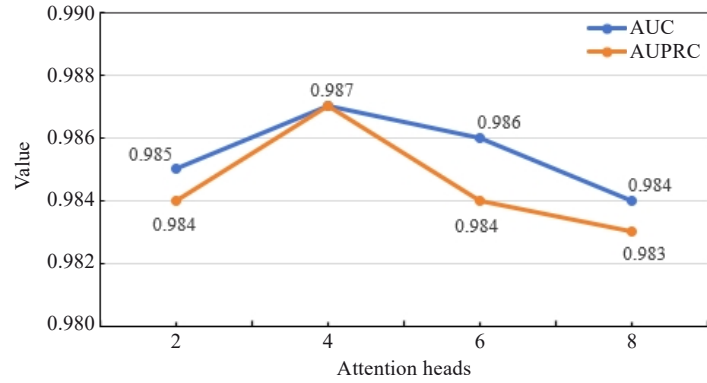


图4 不同注意力头数对GMADTI模型性能的影响
Fig.4 Impact of number of attention heads on the performance of the GMADTI model.

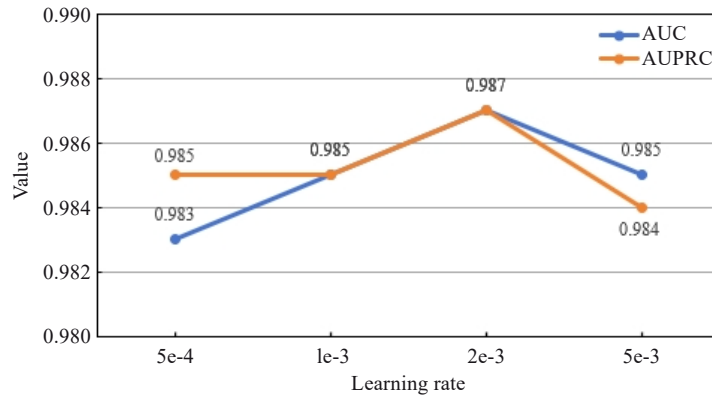


图5 不同学习率对GMADTI模型性能的影响
Fig.5 Impact of learning rates on performance of the GMADTI model.

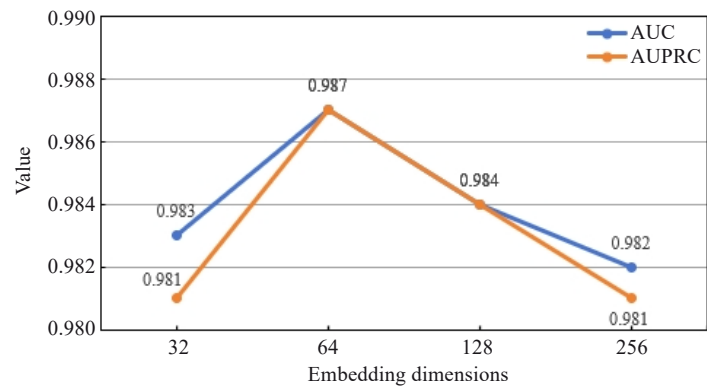


图6 不同嵌入维度对GMADTI模型性能的影响
Fig.6 Impact of embedding dimensions on performance of the GMADTI model.

整体性能、组件贡献、参数敏感性及可视觉解释等方面展开深入分析。首先,从整体性能来看,GMADTI在药物-靶标相互作用预测任务中表现出优异性能。在AUC与AUPRC两项指标上,GMADTI均超越了当前主流的12种基线方法,平均提升到3.2%和3.1%,验证了其在建模复杂异构生物网络方面的有效性。具体而言,在Luo数据集上,GMADTI较性能最佳的AMGDTI提升

了1.0%,展现了更强的药物-靶标关系建模能力。而在更具挑战性的Zheng数据集上,尽管该数据集包含更丰富的药物和蛋白质属性信息,GMADTI依然取得了优异表现,相较于AMGDTI分别提升了1.4%和1.2%,这表明,GMADTI在处理复杂且多样化的药物-靶标关系方面具有显著优势,能够深入理解并利用多源异构数据,同时适应多样化数据集的挑战。

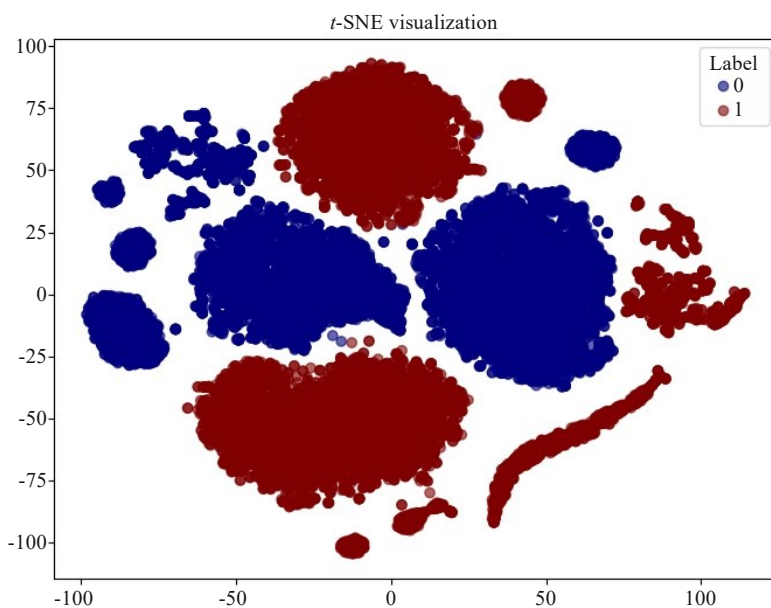


图7 药物与靶标节点的 t-SNE 嵌入可视化分布图
Fig.7 t-SNE visualization of drug and target node embeddings.

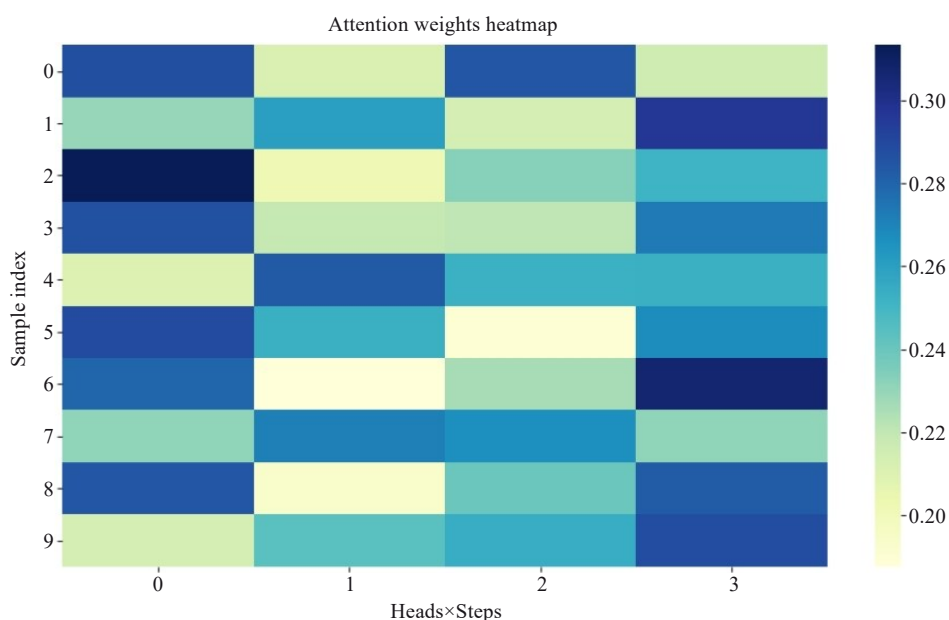


图8 药物侧不同样本的注意力权重分布热力图
Fig.8 Attention weight heatmap for different samples on the drug side.

其次,消融实验进一步验证了模型关键模块的实际贡献。门控机制中动态权重调节对于提升节点特征表达能力和信息传播选择性发挥了关键作用,对模型性能影响显著;多头注意力机制有效增强了语义信息融合的灵活性和上下文适应能力;双重损失函数提升了模型对复杂交互关系的刻画能力;自适应元路径搜索则促进了模型对多样化语义组合的感知与高阶结构信息的挖掘。综合来看,门控机制对整体性能提升贡献最大,是模型设计中的核心组件,其次,多头注意力和双重损失函数在提升模型的语义捕获能力和泛化性能中发挥了重要辅助作用;自适应元路径搜索则进一步增强了模型对多

样化语义组合的感知和高阶结构信息的挖掘能力,对模型性能也有积极贡献。

参数敏感性分析表明,模型对关键超参数如注意力头数、学习率和嵌入维度表现出较好的鲁棒性。具体而言,适量的注意力头数有助于平衡表达能力和计算复杂度,较少的头数(如2)限制了模型对异质图中多样语义特征的并行建模能力,导致特征表达能力不足;而较多的头数(如6或8)会使注意力分布趋于稀疏且分散,增加计算复杂度。在学习率方面,合理的学习率设置对于模型训练的稳定性和收敛速度至关重要。较大的学习率(如5e-3)容易造成参数更新过快,权重波动较大,出

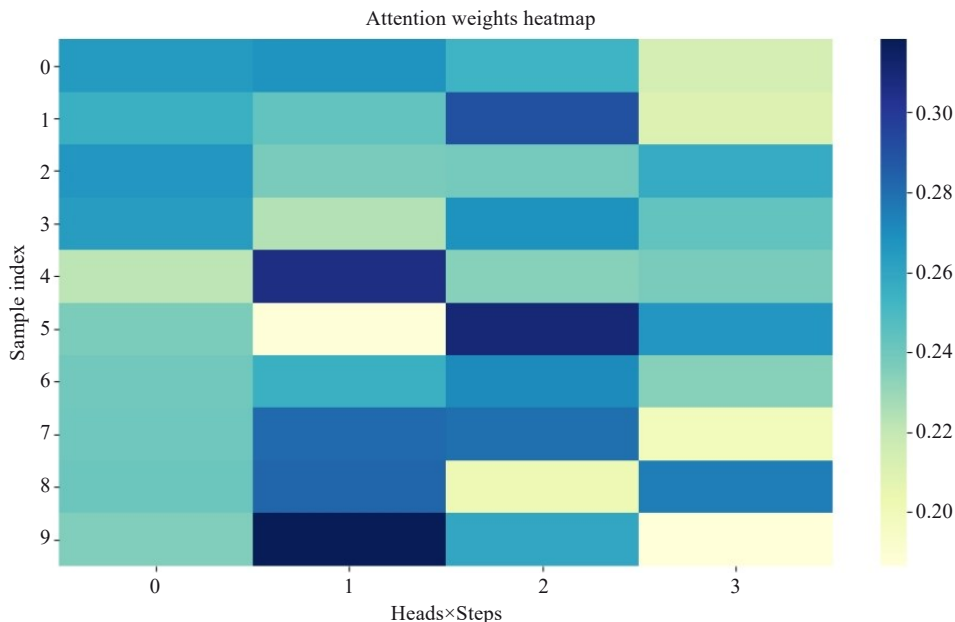


图9 靶标侧不同样本的注意力权重分布热力图

Fig.9 Attention weight heatmap for different samples on the target side.

现过拟合问题;而较小的学习率(如 $5e-4$)则收敛速度缓慢,不利于在有限训练时间内获得理想性能。此外,嵌入维度对模型的表示能力也具有明显影响。当嵌入维度较小时(如32),模型的特征表达能力受限,难以捕捉复杂的关系与语义细节;而当嵌入维度过大(如128或256)时,会引入更多的冗余信息,导致计算成本上升。综上所述,适当的注意力头数、学习率与嵌入维度组合能够显著提升模型对药物-靶标复杂关系的建模能力,从而提高药物-靶标互作预测的准确性与稳定性。

在模型可视化解释方面,t-SNE可视化结果表明,模型在特征空间中能够有效区分药物与靶标节点的语义聚类,节点分布清晰且边界明确,反映出模型在多语义融合和特征提取方面的有效性。同时,注意力热力图揭示了模型在药物视角与靶标视角下对不同语义路径的关注程度存在显著差异,验证了所提出融合多头注意力的语义聚合模块能够自适应整合多层语义信息,并准确评估不同路径的重要性。该机制提升了模型从多个语义层面理解药物-靶标复杂关系的能力,具备较强可解释性。

尽管GMADTI在语义建模、特征表达与预测性能等方面均表现出良好效果,其计算效率仍有进一步提升空间。未来可考虑引入特征压缩策略,以提升模型在大规模图数据上的适用性与扩展性。综上所述,GMADTI通过有效融合多源信息和高阶语义结构,实现了药物-靶标相互作用预测的性能提升,展现出良好的实际应用潜力。未来研究将聚焦于模型的效率优化、泛化能力提升及多场景适应,推动DTI预测向更智能化和实用化方向发展。

本文针对现有药物-靶标相互作用预测方法在高阶

语义建模、自适应语义融合与节点特征过平滑等方面的不足,提出一种基于多层语义与拓扑融合的异质图预测方法GMADTI。该方法通过引入自适应元路径搜索机制,显著增强了模型对高阶语义依赖的捕捉能力;融合多头注意力的语义聚合模块,实现了不同语义路径信息的自适应加权整合;结合结构感知的门控图卷积网络,有效缓解了节点特征冗余与过平滑问题。这些设计共同提升了模型在复杂异质图中的交互关系建模能力。实验结果表明,GMADTI在多个公开数据集上均显著优于现有主流方法,在预测准确性与稳定性方面表现突出。所提出的模型为高维复杂生物网络中信息交互建模提供了新的技术路径,也为药物发现与精准医疗中的理论计算和方法研究提供了有效支撑。未来工作将从以下2个方向开展:结合更多类型的生物分子数据,以及其他生物学信息(如基因组数据、蛋白质结构等),构建大规模异质生物网络模型。优化模型的计算效率,探索更高效的训练方法,如分布式计算。

Declaration of interests: The authors declare no competing interests.

参考文献:

- [1] Abbasi Mesrabadi H, Faez K, Pirgazi J. Drug-target interaction prediction based on protein features, using wrapper feature selection [J]. *Sci Rep*, 2023, 13(1): 3594.
- [2] Dehghan A, Abbasi K, Razzaghi P, et al. CCL-DTI: contributing the contrastive loss in drug-target interaction prediction[J]. *BMC Bioinformatics*, 2024, 25(1): 48.
- [3] 刘正美,魏雪梅,张俊鹏,等. 药物分子与靶标蛋白结合亲和力和预测研究进展[J]. *计算机工程与应用*, 2024, 60(23): 79-90.
- [4] Khatun MS, Hasan MM, Kurata H. PreAIP: computational

- prediction of anti-inflammatory peptides by integrating multiple complementary features[J]. *Front Genet*, 2019, 10: 129.
- [5] Ru XQ, Wang LD, Li LH, et al. Exploration of the correlation between GPCRs and drugs based on a learning to rank algorithm[J]. *Comput Biol Med*, 2020, 119: 103660.
- [6] Ding YJ, Tang JJ, Guo F. Identification of drug-target interactions via multi-view graph regularized link propagation model[J]. *Neurocomputing*, 2021, 461: 618-31.
- [7] Song W, Xu L, Han C, et al. Drug-target interaction predictions with multi-view similarity network fusion strategy and deep interactive attention mechanism[J]. *Bioinformatics*, 2024, 40(6): btae346.
- [8] Luo Y, Zhao X, Zhou J, et al. A network integration approach for drug-target interaction prediction and computational drug repositioning from heterogeneous information[J]. *Nat Commun*, 2017, 8(1): 573.
- [9] Zhou D, Xu Z, Li W, et al. MultiDTI: drug-target interaction prediction based on multi-modal representation learning to bridge the gap between new chemical entities and known heterogeneous network[J]. *Bioinformatics*, 2021, 37(23): 4485-92.
- [10] Qiao G, Wang G, Li Y. Causal enhanced drug-target interaction prediction based on graph generation and multi-source information fusion[J]. *Bioinformatics*, 2024, 40(10): btae570.
- [11] 秦海盈, 赵中英, 李建晖, 等. 基于元路径与层次注意力的时序异质信息网络表示学习方法[J]. *模式识别与人工智能*, 2021, 34(12): 1093-102.
- [12] Li J, Wang J, Lv H, et al. IMCHGAN: inductive matrix completion with heterogeneous graph attention networks for drug-target interactions prediction[J]. *IEEE/ACM Trans Comput Biol Bioinform*, 2022, 19(2): 655-65.
- [13] Su Y, Hu Z, Wang F, et al. AMGDTI: drug-target interaction prediction based on adaptive meta-graph learning in heterogeneous network[J]. *Brief Bioinform*, 2023, 25(1): bbad474.
- [14] 潘柏儒, 丁卫平, 鞠恒荣, 等. 基于粗糙图的图卷积神经网络算法[J]. *模式识别与人工智能*, 2023, 35(9): 827-38.
- [15] Wan F, Hong L, Xiao A, et al. NeoDTI: neural integration of neighbor information from a heterogeneous network for discovering new drug-target interactions[J]. *Bioinformatics*, 2019, 35(1): 104-11.
- [16] Zhao T, Hu Y, Valsdottir LR, et al. Identifying drug-target interactions based on graph convolutional network and deep neural network[J]. *Brief Bioinform*, 2021, 22(2): 2141-50.
- [17] Peng J, Wang Y, Guan J, et al. An end-to-end heterogeneous graph representation learning-based framework for drug-target interaction prediction[J]. *Brief Bioinform*, 2021, 22(5): bbaa430.
- [18] Zhu Y, Ning C, Zhang N, et al. GSRF-DTI: a framework for drug-target interaction prediction based on a drug-target pair network and representation learning on a large graph[J]. *BMC Biol*, 2024, 22(1): 156.
- [19] Mukesh M, Raghuvanshi Pradnya S, Borkar Sachin U, Balvir. Node2Vec and machine learning: a powerful Duo for link prediction in social network[J]. *Jes*, 2024, 20(2s): 639-49.
- [20] Johnson SJ, Murty MR, Navakanth I. A detailed review on word embedding techniques with emphasis on word2vec[J]. *Multimed Tools Appl*, 2024, 83(13): 37979-8007.
- [21] Ding Y, Yao Q, Zhao H, et al. Diffing: Differentiable meta graph search for heterogeneous graph neural networks[C]. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021: 279-88.
- [22] Guo YB, Zhou DM, Ruan XL, et al. Variational gated autoencoder-based feature extraction model for inferring disease-miRNA associations based on multiview features[J]. *Neural Netw*, 2023, 165: 491-505.
- [23] Rasamoelina AD, Adjailia F, Sincak P. A review of activation function for artificial neural network[C]//2020 IEEE 18th World Symposium on Applied Machine Intelligence and Informatics (SAMII). January 23-25, 2020. Herlany, Slovakia. IEEE, 2020: 281-6.
- [24] Dubey SR, Singh SK, Chaudhuri BB. Activation functions in deep learning: a comprehensive survey and benchmark[J]. *Neurocomputing*, 2022, 503: 92-108.
- [25] Deora P, Ghaderi R, Taheri H, et al. On the optimization and generalization of multi-head attention[J]. *arXiv preprint arXiv: 2310.12680*, 2023.
- [26] 郑子瑜, 杨夏颖, 吴圣杰, 等. 多特征融合的产时超声胎方位识别模型[J]. *南方医科大学学报*, 2025, 45(7): 1563-70.
- [27] Tian Y, Wang X, Yao X, et al. Predicting molecular properties based on the interpretable graph neural network with multistep focus mechanism[J]. *Brief Bioinform*, 2023, 24(1): bbac534.
- [28] Wang F, Liu H. Understanding the behaviour of contrastive loss[C]. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021: 2495-504.
- [29] Ranjan P, Khan P, Kumar S, et al. Slog-Sigmoid activation-based long short-term memory for time-series data classification[J]. *IEEE Trans Artif Intell*, 2024, 5(2): 672-83.
- [30] Zheng Y, Peng H, Zhang XC, et al. Predicting drug targets from heterogeneous spaces using anchor graph hashing and ensemble learning[C]//2018 International Joint Conference on Neural Networks (IJCNN). July 8-13, 2018. Rio de Janeiro. IEEE, 2018: 1-7.
- [31] Srivastava A, Rawat BS, Singh G, et al. A review of optimization algorithms for training neural networks[C]//2023 International Conference on Sustainable Emerging Innovations in Engineering and Technology (ICSEIET). September 14-15, 2023. Ghaziabad, India. IEEE, 2023: 886-90.
- [32] Wang H, Huang F, Xiong Z, et al. A heterogeneous network-based method with attentive meta-path extraction for predicting drug-target interactions[J]. *Brief Bioinform*, 2022, 23(4): bbac184.
- [33] Li Y, Qiao G, Gao X, et al. Supervised graph co-contrastive learning for drug-target interaction prediction[J]. *Bioinformatics*, 2022, 38(10): 2847-54.
- [34] Yao KN, Wang XW, Li WN, et al. Semi-supervised heterogeneous graph contrastive learning for drug-target interaction prediction[J]. *Comput Biol Med*, 2023, 163: 107199.

(编辑:林萍)